# Instrument Response Studies

## Agenda

- Overarching Approach & Strategy
- Classification Trees
- Sorting out Energies
- PSF Analysis
- Background Rejection
- Assessment

GLAST

# Overarching Approach & Strategy

## A 3 Stage Approach

1. Energy determination   -   <span style="color:green">Foundational to what follows</span>

2. Evaluate PSF's   -   <span style="color:blue">Background will be suppressed</span>

3. Reject the Background   -   <span style="color:red">The hard part</span>

Statistical Tools:  Classification Trees & Regression Trees

# A Brief History of Resolution & Rejection

## Preparing for DC1 is a LARGE TASK

- Not likely to get right the 1st, or the 2nd, or the 3rd, or.... time!

1st Time: April-May
> Discover Mult-scattering in G4 "too good to believe!"
> Took till end of June to fix!

2nd Time: July (SAS Workshop)
> OOPS!  The ACD geometry!

3rd Time: July-August
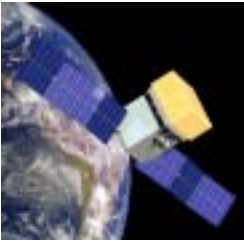> Where did all the Run Numbers go?

4th Time: August
> Will Bill never stop changing variable - well at least
> he shouldn't make so many coding errors! Steve's variables added.

5th Time: August-September
> Data of the day!  But its certainly not "The rest of the story!"

6th Time: .... IS A CHARM!

GLAST

# Classification Tree Primer

Origin:  Social Sciences  - 1963
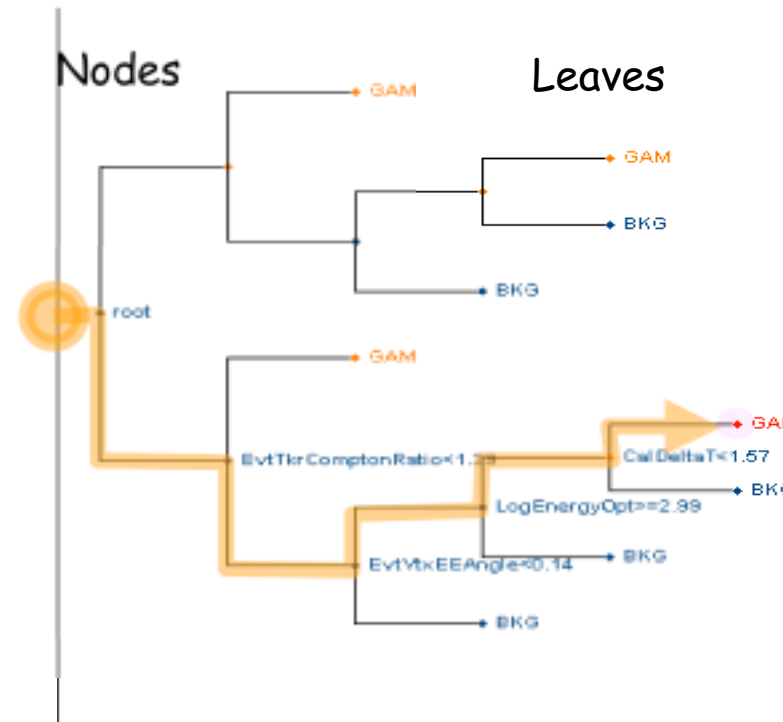
How a CT works is simple:
    A series of  "cuts" parse the
    data into a "tree" like structure,
    where final nodes (leaves) are "pure"

A "traditional analysis" is just ONE path
through such a tree.
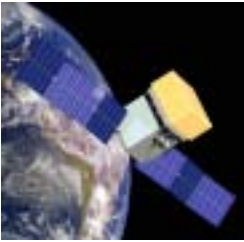
Tree are *much* more efficient!

Mechanism of tree generation less subject
to "investigator basis."

A Simple Classification Tree

Nodes | Leaves

GAM

GAM

BKG

BKG

GAM

root

GAM

EvtTkrComptonRatio<1. | CalDeltaT<1.57

LogEnergyOpt>=2.99 | BKG

EvtVtxEEAngle<0.14 | BKG

BKG

STATISTICALLY HONEST!

GLAST

# Input Data for Training and Testing

"Tree Production" automated by using "Training Samples" where the results are *a priori* known

All-Gammas (AG):   $18 \text{ MeV} < E_\gamma < 18 \text{ GeV}$

  $1/E$ Spectrum

  $-1 < \cos(\theta) < 0$ (2π str)

  $A_{GEN} = 6 \text{ m}^2$

AG Total: $3/4 \times 10^6$ Events

CAL -Training         25%
PSF -Training         50%
BKG -Training/Testing 25%

Background Events(BGEs):   0:  Orbit Ave CHIME
  1:  Albedo Protons
  2:  Albedo $\gamma$s
  3:  Cosmic $e^-$
  4:  Albedo $e^+$ & $e^-$

  $A_{GEN} = 6 \text{ m}^2_5$

BKG Total: $.9 \times 50 \times 10^6$ Events

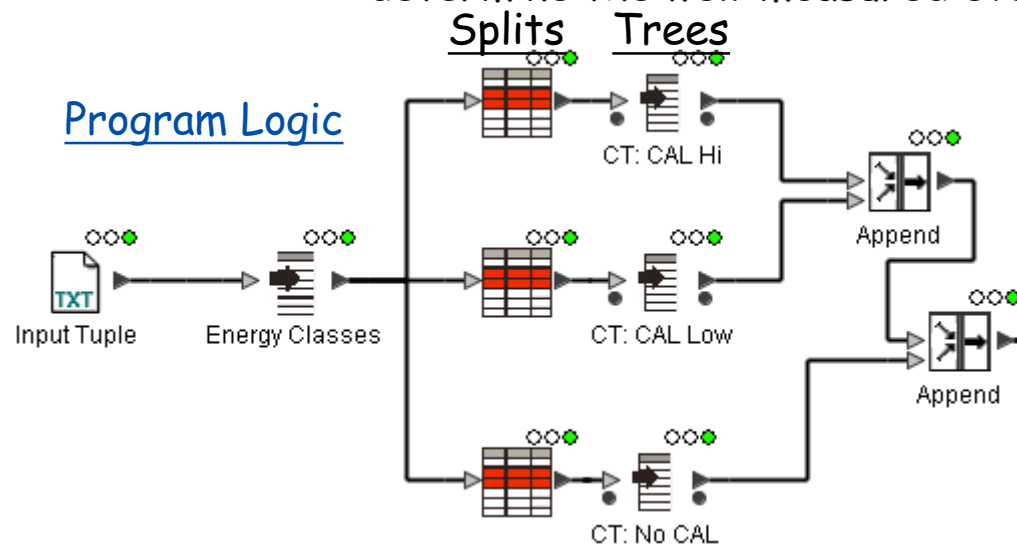BKG -Training  50%
BKG -Testing   50%

GLAST

# Energy Filtering

Problem: The large gaps in the CAL and the thick layers of the Tracker
compromise the energy determination.

Strategy: Identify poorly measured events and eliminate them.

Technique: Split events into classes and for each class use a Classification Tree to
determine the well-measured events.



Energy Class Definitions

CAL-Hi:  CalEnergySum > 100 MeV
        CalTotRLn > 2

CAL-Low: CalEnergySum < 100 MeV
        CalEnergySum > 5 MeV
        CalTotRLn > 2

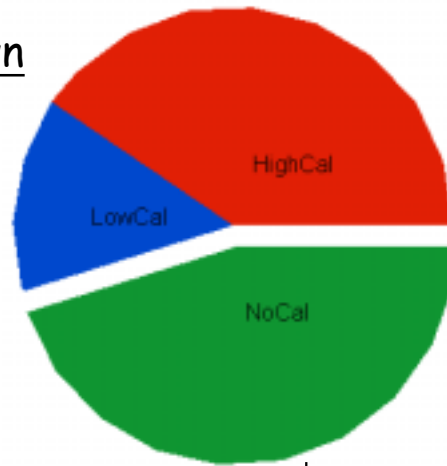No-CAL:  CalEnergySum < 5 MeV or
        CalTotRLn < 2

GLAST

# Energy Filtering (2)

Energy Class Breakdown

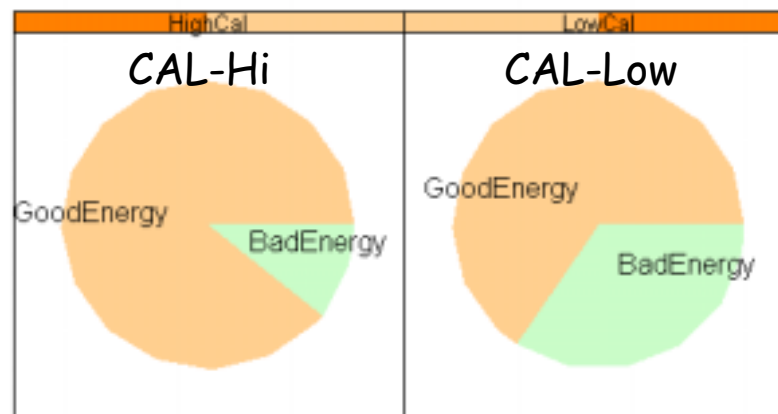CAL-Hi: 41%

CAL-Low: 14%

No-CAL: 45%

The No-CAL are presently not analyzed.

These will need to be addressed in the future as it constitutes the largest Energy Event Class and could greatly improve the transient response

CT Energy Classes: "GoodEnergy" = $\left(\sigma_{Energy} < 35\%\right)$ $\left|\dfrac{CalEnergySumOpt - McEnergy}{McEnergy}\right| < .35$

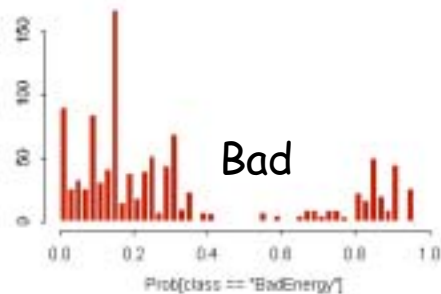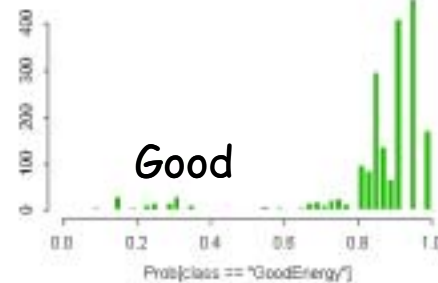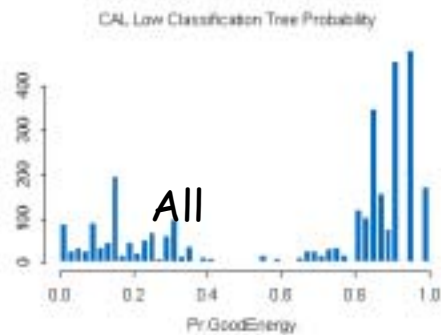"GoodEnergy" / "BadEnergy"
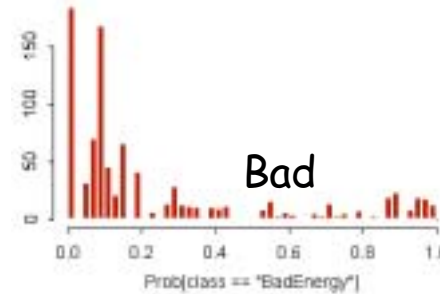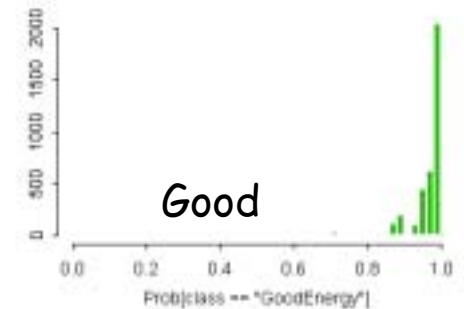
Event Breakdown by Energy Class

GLAST

# Energy Filtering (3)

All available variables bearing on the quality of the
energy determination are made available to "train"

### CAL-Low CT Probabilities



### CAL-High CT Probabilities

# Energy Filtering (4)

Cut:
Cal.Prob > .50

Eff. = 82%

Bad-Cal = 4.5%

Before

After

GLAST

# Energy Filtering (5)

The Results:

Cut more severe as events near Instrument Axis

We can use this for SCIENCE!



Over Estimates "Clean"

Some Low Energy Straglers

EvtLogESum

McLogEnergy

# PSF Filtering

**Program Logic**



Global Cuts:
1) Cal.Prob > .50                    (-18%)

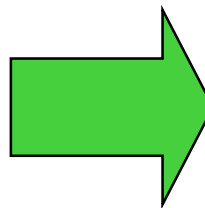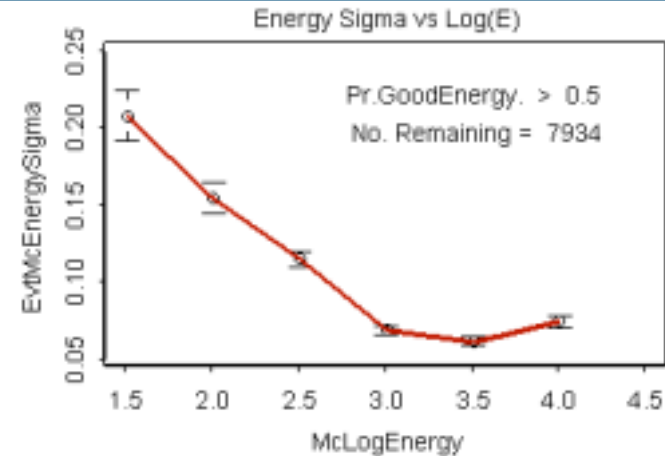Cleaning Cuts Applied to CT Training
2) EvtTkr1EChisq < 7.5          &
   EvtTkr1EFirstChisq < 10.  &
   EvtTkr2EChisq < 10.          &
   EvtTkr2EFirstChisq < 10  (-5.6%)

TOTAL LOSS:   -22.5% (Training)
                      -18%    (Analysis)

Thin / Thick Split:   Best Track originates in Thin / Thick Radiators
                                48% Thin  / 52% Thick

VTX / 1Tkr Split:  Use CT to determine whether or not to use Recon VTX Solution

1 CT & 1 RT Used for each of the 4 PSF Classes:  CT used to kill long tail
                                RT used to sharpen CORE resolution

GLAST

# PSF Filtering: VTX/1Tkr Split

Only events with a VTX solution are considered (VtxAngle > 0)

Using MC Truth, the best solution is determined (for CT Training)

Mariginal Improvement:
Purity (Before/After)  60% / 66%
(See Discussion at end of talk)



**Input Node - Filter Rows (1084)**

| | | Predicted | | Totals |
|---|---|---|---|---|
| | | 1TKR | VTX | |
| Observed | 1TKR | 399 | 555 | 954 |
| | VTX | 406 | 1052 | 1458 |
| Totals | | 805 | 1607 | 2412 |

| | Observed | | Overall |
|---|---|---|---|
| | 1TKR | VTX | |
| % Agree | 41.8% | 72.2% | 60.2% |

**Positive Category - VTX**

| Recall | Precision | F-Measure |
|---|---|---|
| 72.2% | 65.5% | 68.6% |



Relative Column Importance

GLAST

# PSF Tails

"Tail" Events defined as being 2.3 x *PSF Model* or worse.

Improvement:

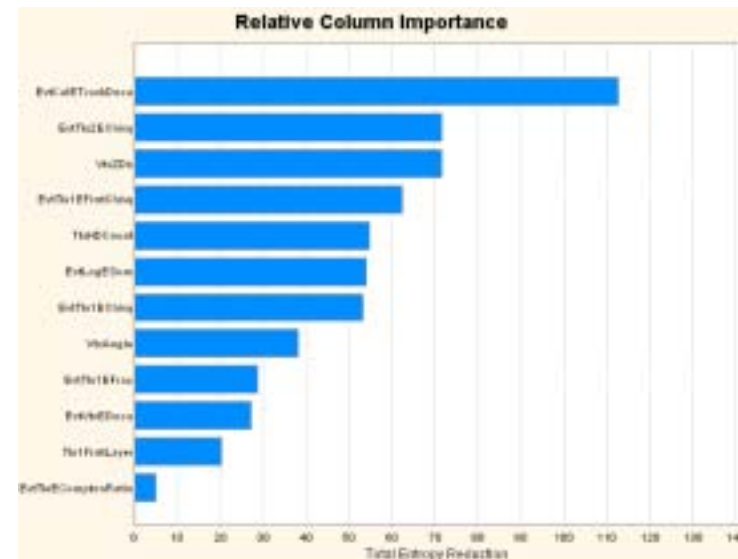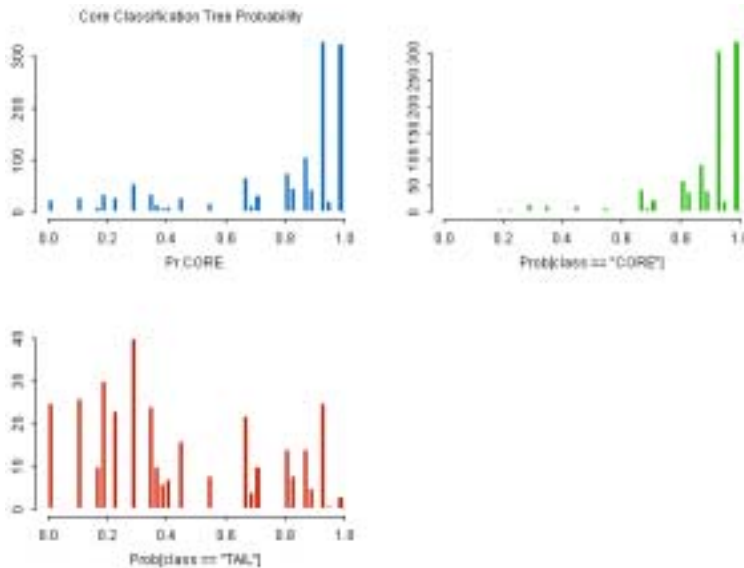38% of the "Tail" is eliminated at expense of 13.5% of the "Core"

**Input Node – Predict: VTX Class (1095)**

| | | Predicted | | Totals |
|---|---|---|---|---|
| | | CORE | TAIL | |
| Observed | CORE | 1244 | 195 | 1439 |
| | TAIL | 104 | 64 | 168 |
| Totals | | 1348 | 259 | 1607 |

| | Observed | | Overall |
|---|---|---|---|
| | CORE | TAIL | |
| % Agree | 86.4% | 38.1% | 81.4% |

**Positive Category – CORE**

| Recall | Precision | F-Measure |
|---|---|---|
| 86.4% | 92.3% | 89.3% |

GLAST

# PSF CORE

Tool: Regression Tree (Similar to CT)
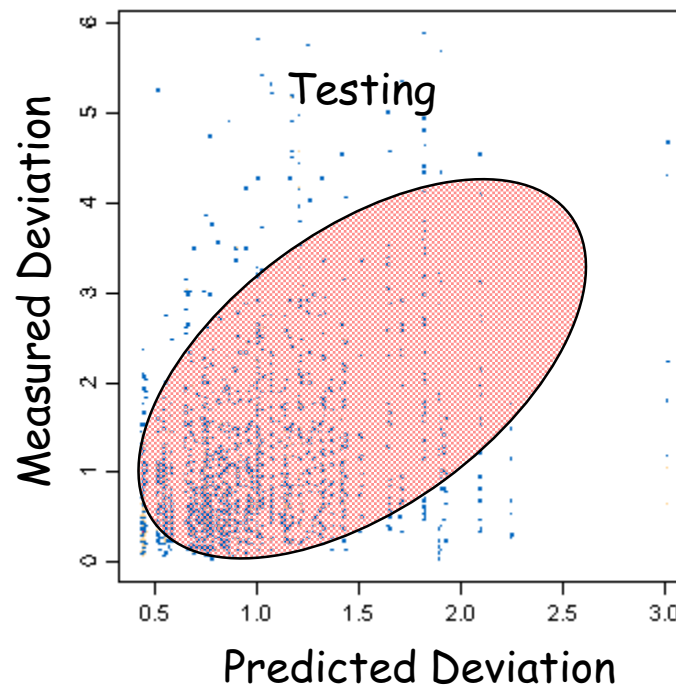   Matches deviations rather then
   class types.
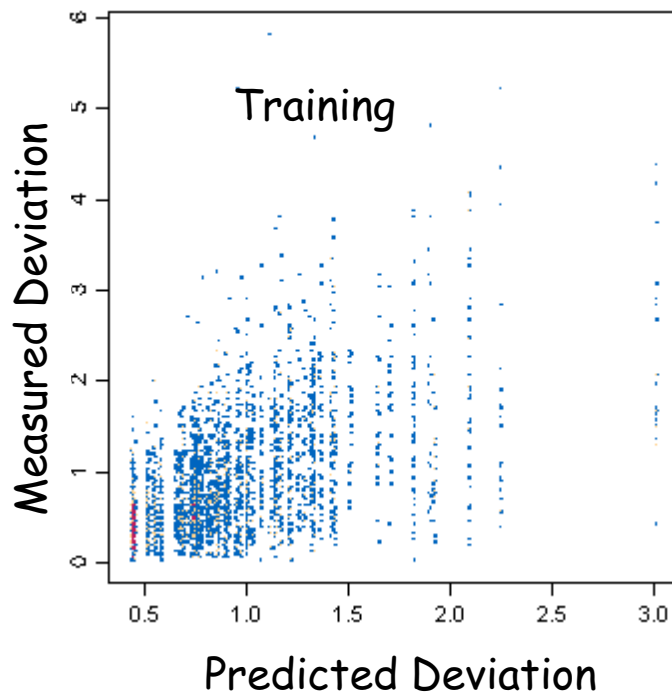
Event-by-Event PSF Error

Energy Compensated by: $1/E_{Meas}^{.8}$

Collapse All PSF's onto one.
Normalization:  1 = PSF(68) Sci. Req.

Event Starvation VERY APPARENT!



Training — Measured Deviation vs. Predicted Deviation
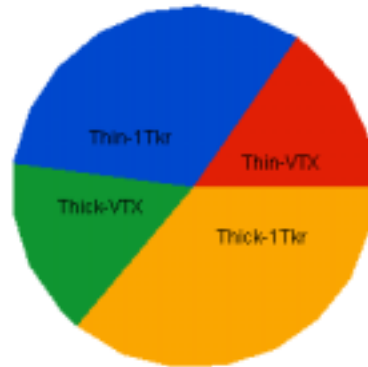
Testing — Measured Deviation vs. Predicted Deviation
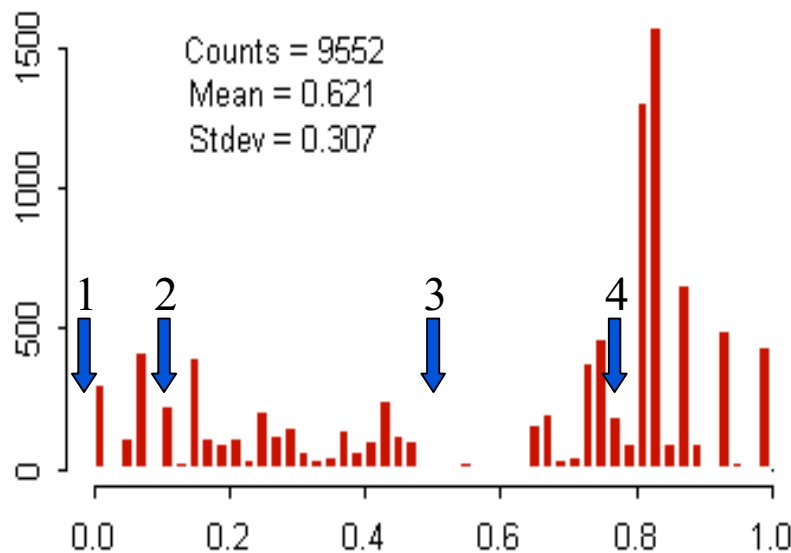
GLAST

# PSF Summary

PSF Class Breakdown:
Thin-VTX: 15.3%
Thin-1Tkr: 32.7%
Thick-VTX: 15.9%
Thick-1Tkr: 36.0%

PSF Clean-up Cuts:

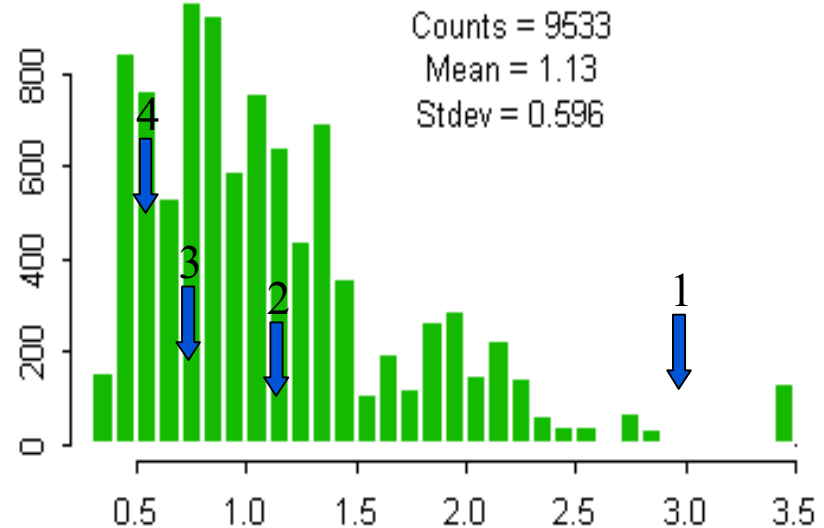Matrix of 4x4 PSF Plots vs Log(E) examined



PSF Probability Distributions

Counts = 9552
Mean = 0.621
Stdev = 0.307

Core Cut: Limit PSF tails

PSF Probability Distributions
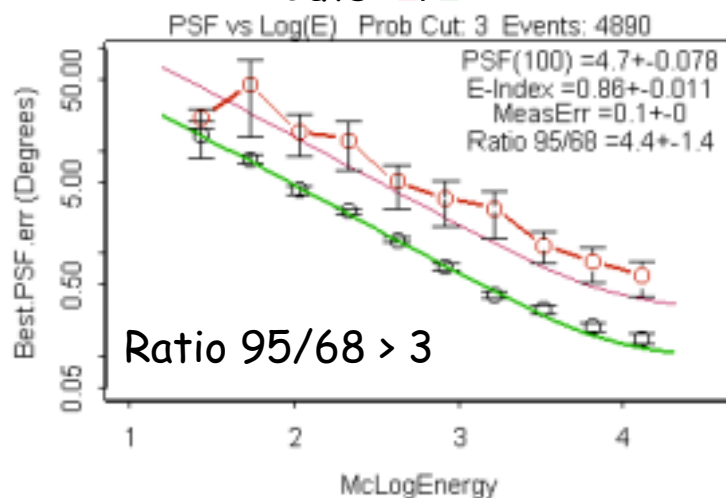
Counts = 9533
Mean = 1.13
Stdev = 0.596

Pred. PSF: Sharpen PSF

# Thin PSF's - Integrated over FoV
## 4 Combinations of Cuts *(CORE/Pred)*

### Cuts: **1**/**1**

PSF vs Log(E)   Prob Cut: 3   Events: 4890

PSF(100) =4.7+-0.078
E-Index =0.86+-0.011
MeasErr =0.1+-0
Ratio 95/68 =4.4+-1.4

Best.PSF.err (Degrees)

McLogEnergy

Ratio 95/68 > 3

### Cuts: **2**/**1**

PSF vs Log(E)   Prob Cut: 3   Events: 4621

PSF(100) =4.1+-0.28
E-Index =0.8+-0.044
MeasErr =0.1+-0
Ratio 95/68 =3+-0.7

Best.PSF.err (Degrees)

McLogEnergy

Meets SR
Events Eff.: 94.5%

### Cuts: **3**/**2**

PSF vs Log(E)   Prob Cut: 1.1   Events: 2557

PSF(100) =3.6+-0.11
E-Index =0.86+-0.019
MeasErr =0.1+-0
Ratio 95/68 =2.3+-0.2

Best.PSF.err (Degrees)

McLogEnergy

Events Eff.: 52.3%

### Cuts: **3**/**4**

PSF vs Log(E)   Prob Cut: 0.55   Events: 935

PSF(100) =3+-0.22
E-Index =0.82+-0.055
MeasErr =0.1+-0
Ratio 95/68 =2.2+-0.33

Best.PSF.err (Degrees)

McLogEnergy

Events Eff.: 19.1%

GLAST
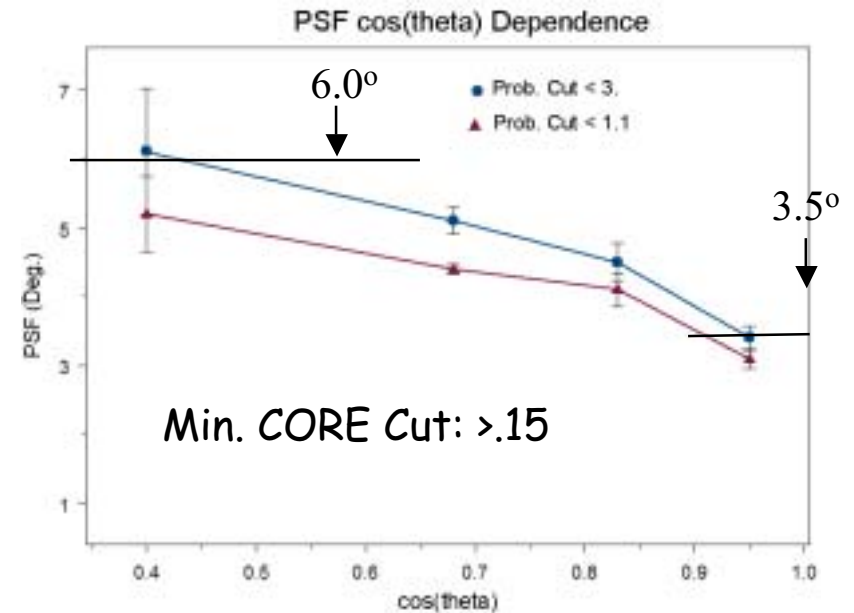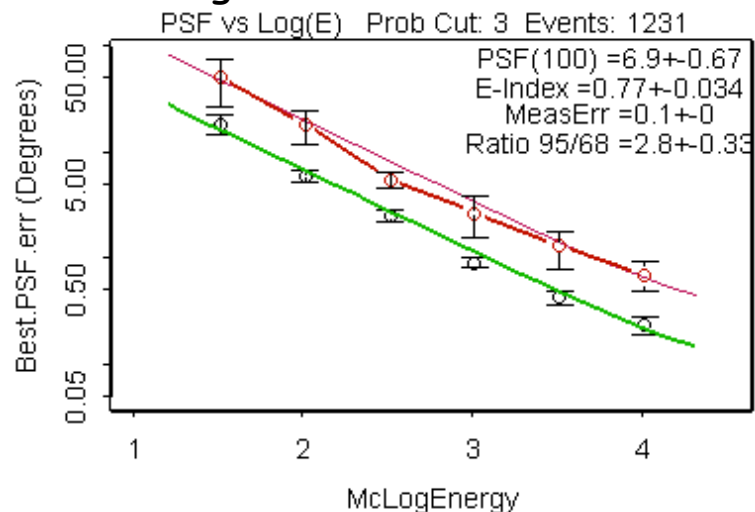
# PSF Summary - Minimum CORE Cut

PSFs given prior to Background Rejection due to lack of statistics

Background rejection does not change conclusions.

Limited statistics don't allow for good determination of PSF vs $\cos(\theta)$ for tight cuts



PSF cos(theta) Dependence

6.0°

3.5°

● Prob. Cut < 3.

▲ Prob. Cut < 1.1

Min. CORE Cut: >.15



PSF vs Log(E)   Prob Cut: 3  Events: 1231

PSF(100) =6.9+-0.67
E-Index =0.77+-0.034
MeasErr =0.1+-0
Ratio 95/68 =2.8+-0.33

## Thick Radiator PSF

PSF(Thick) = 2 x PSF(Thin)

CORE Cut and Pred. CORE are adjusted to have similar effects as for Thin Radiators

# $A_{eff}$ Summary - Minimum CORE Cut

Lack of events makes determination imprecise!

Effective Area On Axis ($E_\gamma > 3$ GeV)

$A_{eff} = N_{Obs}/N_{Gen}$ x 6 x 1.3

$A_{eff} = 2603/18750$ x 7.8

$A_{eff} = 1.1$ m$^2$

Light Gathering Power ($E_\gamma > 3$ GeV)

$A_{eff}$ x $\Delta\Omega = N_{Obs}/N_{Gen}$ x 6 x 2$\pi$ x 1.27

$A_{eff}$ x $\Delta\Omega = 9877/187500$ x 37.7 x 1.27

$A_{eff}$ x $\Delta\Omega = 2.5$ m$^2$-str

Angular Dependence

~ Linear in cos($\theta$)

At low energy FoV is truncated

Slight roll-over near axis due to CAL inefficiency caused by inter-tower gaps



Effective Area - SR Cuts

135/bin asymptotic

Effective Area - SR Cuts

Note:
On Axis Roll-Off
-.80<cos($\theta$)<-.60

-1<cos($\theta$)<-.80

-.40<cos($\theta$)<-.20

-.60<cos($\theta$)<-.40

Tkr1ZDir

# Background Rejection

## Pre-Analysis Filtering
Done to reduce data volume

Require at least 1 Reconstructed Track

Require AcdActiveDist < -20 mm
(AcdActiveDist defined to be distance to
 edge of nearest hit Acd Tile.  Values < 0
 indicate projected track falls
 OUTSIDE of hit tile area.)
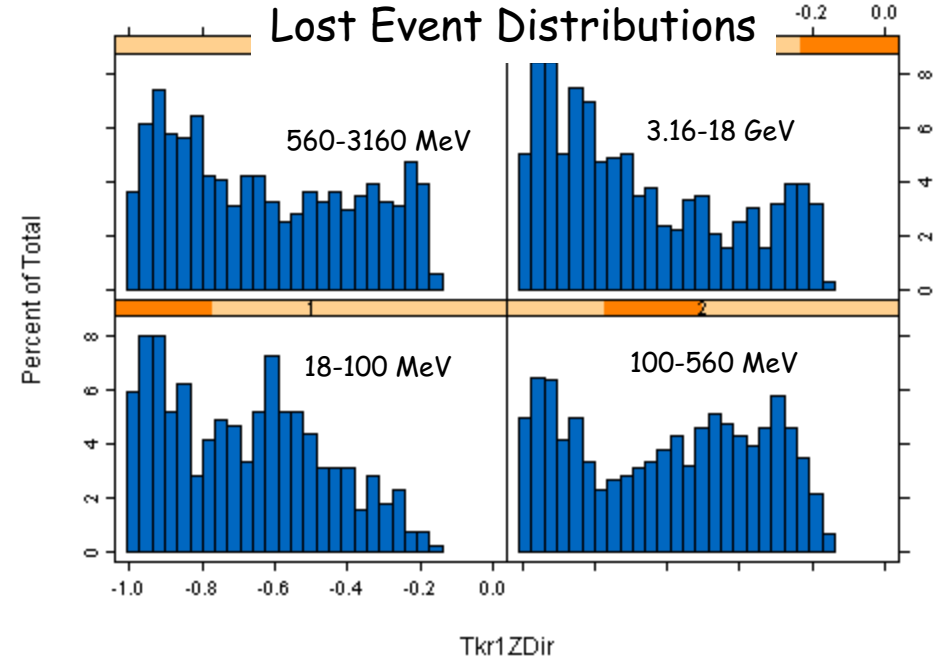
Note: This has a built in Energy Dependence!

| | |
|---|---|
| Generated: | $50 \times 10^6$ |
| Lost 10% from failed jobs: | $45 \times 10^6$ |
| Number of Triggers: | $\sim 18.5 \times 10^6$ |
| Number left after pre-filter: | $.73 \times 10^6$ |

## First Analysis Cut:
Require "GoodCal" Energy
Results in 18% loss in $\gamma$ Events
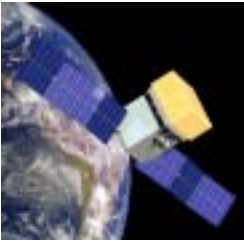Distribution of Event Loss in $\cos(\theta)$

Lost Event Distributions

560-3160 MeV

3.16-18 GeV

18-100 MeV

100-560 MeV

Percent of Total

Tkr1ZDir

Background Event Efficiency: 12.2%
BGE Left: $89.3 \times 10^3$
BGE Trigger Reduction Factor: ~200

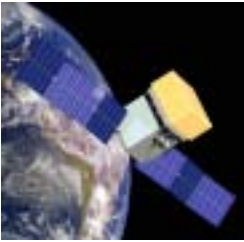# Background Rejection Event Files

BGE sample divided in 2:

    50% Training for CT's

    50% Testing results

    (44652 Events in each)

    Remaining AG sample (25% of original)

        50% Training (12.5% of original)

        50% Testing  (12.5% of original)

      BGE's and AG's tagged and mixed randomly together for both Training and Testing

## This leaves to few events to do much more then explore BGE Rejection problem areas.
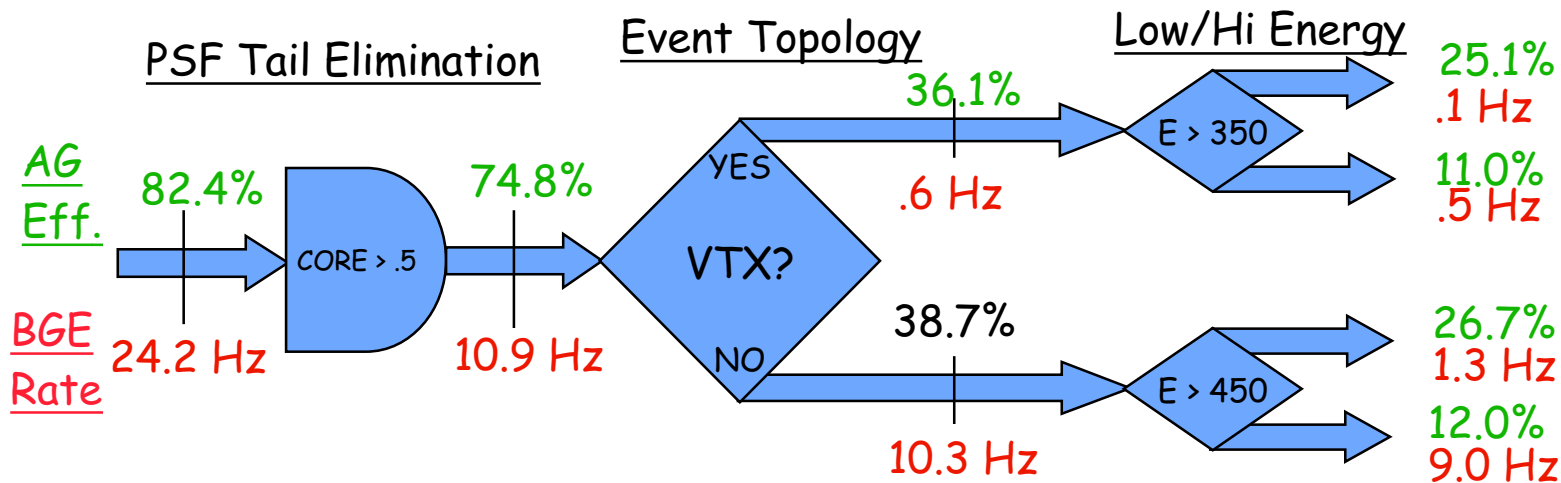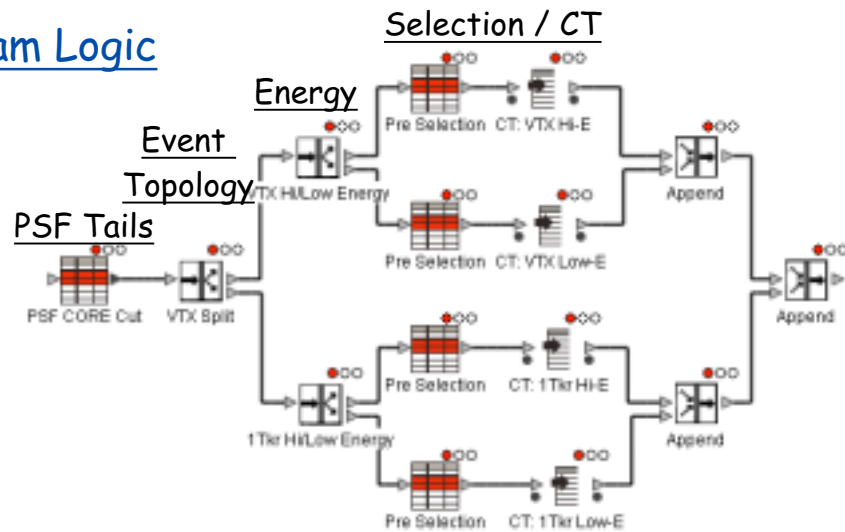
    (i.e. 5629 AG's in each)

GLAST

# Background Rejection Program

Events with a found VTX have much less background

Large energy dependence suggests subdividing into Low/Hi branches
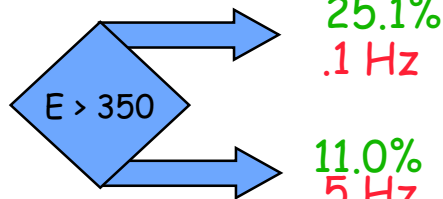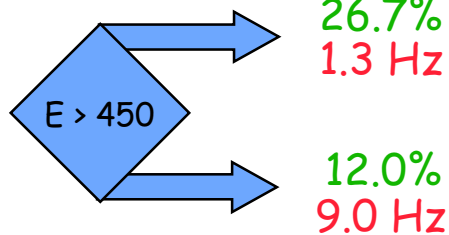
Large rejection Variables used in Pre Selections

Program Logic

GLAST

# Background Rejection Program - Pre Selection

### Pre Selection Cuts

**Low/Hi Energy**

$$E > 350$$

25.1%
.1 Hz

11.0%
.5 Hz

**AG**
**Eff.**

**BGE**
**Rate**

$$E > 450$$

26.7%
1.3 Hz

12.0%
9.0 Hz

| EvtTkrEComptonRatio > .60 & CalMIPDiff > 60. |
|---|
| AcdTileCount == 0 & CalMIPDiff > -125 & EvtTkrEComptonRatio > .80 |

| AcdTotalEnergy < 6.0 & EvtTkrComptonRatio > .70 & CalMIPDiff > 80. & CalLRmsRatio < 20. |
|---|
| AcdTileCount == 0 & EvtTkrComptonRatio > 1. & CalLRmsRatio > 5. & Tkr1FirstLayer != 0 & Tkr1FirstLayer < 15 |

**Out of**

23.2%
.04 Hz

27.4%
(84.7%)

8.4%
.08 Hz

20.7%
(40.6%)

% in Blue show
Rel. Eff. to Event
Sample in that Branch

23.1%
.26 Hz

27.8%
(83.1%)

5.5%
.25 Hz

24.3%
(22.6%)

# Background Rejection Program - CT's

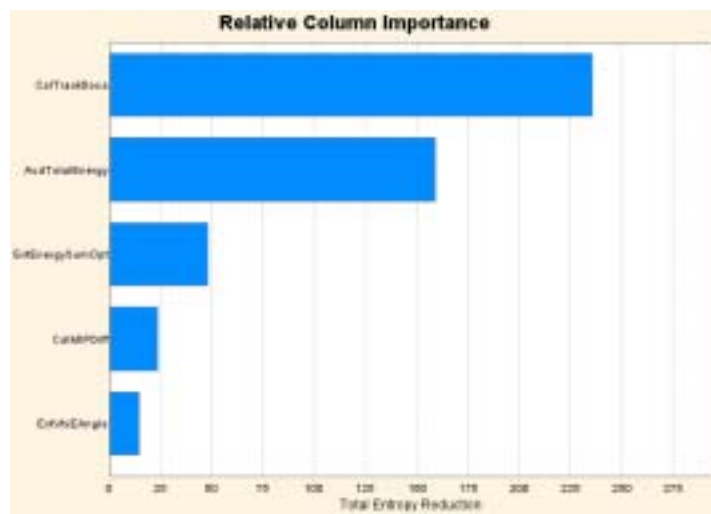**VTX & Hi-E Case**

Training Sample

Note the lack of events!



Input Node - Classification Tree (99)

| | | Predicted | | Totals |
|---|---|---|---|---|
| | | BKG | GAM | |
| Observed | BKG | 81 | 4 | 85 |
| | GAM | 13 | 941 | 954 |
| | Totals | 94 | 945 | 1039 |

| | Observed | | Overall |
|---|---|---|---|
| | BKG | GAM | |
| % Agree | 95.3% | 98.6% | 98.4% |

Positive Category - GAM

| Recall | Precision | F-Measure |
|---|---|---|
| 98.6% | 99.6% | 99.1% |

Background Rejection Tree Probability

Few Events results in sparse CT Trees

**Relative Column Importance**

Total Entropy Reduction

## Testing Results Retention:
AG: 97.5%
BGE: 22.%

Input Node - Predict: Classification Tree (223)

| | | Predicted | | Totals |
|---|---|---|---|---|
| | | GAM | BKG | |
| Observed | GAM | 1544 | 40 | 1584 |
| | BKG | 18 | 65 | 83 |
| | Totals | 1562 | 105 | 1667 |

| | Observed | | Overall |
|---|---|---|---|
| | GAM | BKG | |
| % Agree | 97.5% | 78.3% | 96.5% |

Positive Category - GAM

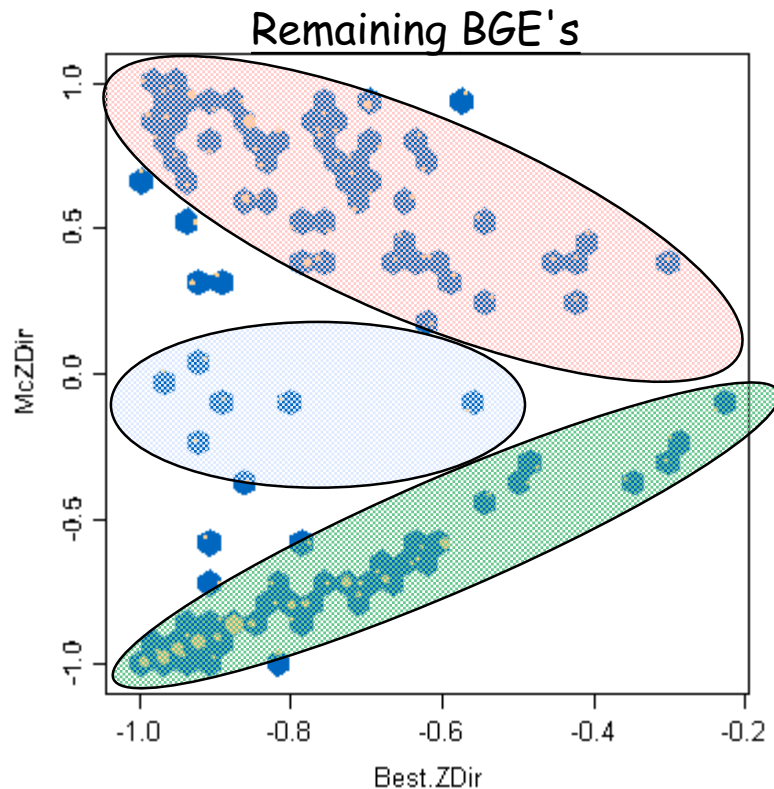| Recall | Precision | F-Measure |
|---|---|---|
| 97.5% | 98.8% | 98.2% |

GLAST

# Background Rejection Program - CT Results

| Case | | CT Tree Disc. | | Out of |
|------|---|---------------|---|--------|

**23.2%**
.04 Hz  →  Hi- E

Prob.Gam > .5

**22.6%**
.01 Hz

27.4%
(82.5%)

VTX (350 MeV)

**8.4%**
.08 Hz  →  Low- E

Prob.Gam > .9

**5.0%**
.02 Hz

20.7%
(24.2%)

**23.1%**
.26 Hz  →  Hi- E

Prob.Gam > .5

**21.5%**
.02 Hz

27.8%
(77.3%)

1Tkr (450 MeV)

**5.5%**
.25 Hz  →  Low- E

Prob.Gam > .9

**1.8%**
.02 Hz

24.3%
(7.4%)

GLAST

# Background Rejection Program - What's Left?

### Remaining BGE's



**3 Classes of BGE Events Remain:**

1) 1:1 Correlated Events - ACD Leakage and inefficiency (.04 Hz)

2) 1: -1 Correlated Events - Range-outs from below (.025 Hz)

3) Events at McZDir ~ 0 - Horizontal Events (.005 Hz)
   <u>Elimination Strategy</u>

1) ACD Leakage
   - Events found accurately;
   - Small phase space
   - Track projection to ACD cracks

2) Range-outs - MIP Identification in CAL

3) Horizontal Events - Edge CAL hits

$A_{eff}$ & BGE Rate:

$A_{eff}$ = 8400 cm$^2$ on Axis (E > 3 GeV)

$A_{eff} \times \Delta\Omega$ = 2.0 m$^2$-str

   BUT....

BGE Rate 5X too high

GLAST

# Back to CT Basics

CT Tree Generation Mechanism:

Variable Selection: $\left|\dfrac{\langle good \rangle - \langle bad \rangle}{\sqrt{\sigma^2_{good} + \sigma^2_{bad}}}\right|$
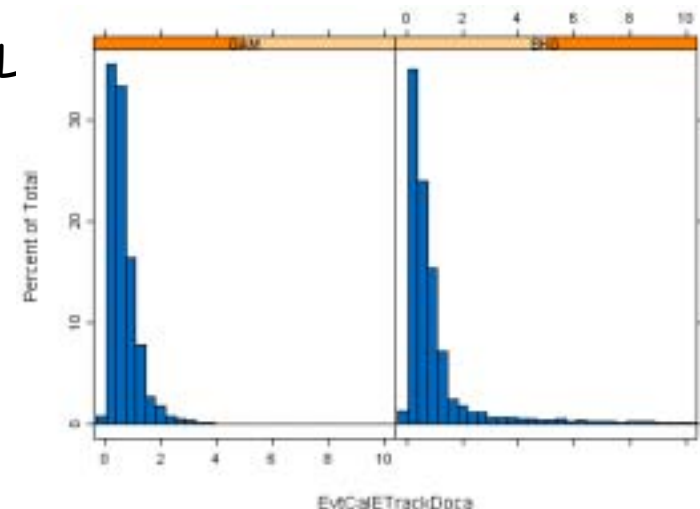
This is a FIRST ORDER TECHNIQUE

When MEANS are approx. equal it fails!

This is the case for MOST OF GL

Example:

One of the most useful
separation variables:
Energy compensated
Cal-Centroid - Track distance

Means similar - Tails dissimilar

GLAST

# A New CT Mechanism

1. Characterize Distribution extents (tails) by <u>Quantiles</u>

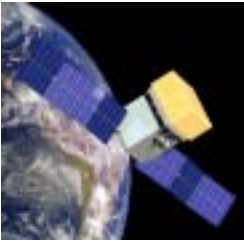   Example: 95% containment PSF is the 95[th] Quantile of the PSF distribution

   Alternative Variable Selection:

   Q(Good, 95) - Q(Bad, 95)  or - normalized... $\left| \dfrac{Q(Good,95) - Q(Bad,95)}{\sqrt{\sigma_{Good} \cdot \sigma_{Bad}}} \right|$

   Use Generic   $N \cdot \log(N)$   for cut placement.

2. CT Generation is a "one step look ahead" - extend to 2,3, etc. steps

3. More Advanced CT Technologies - Ensembles, Boosted Trees, etc.

GLAST

# Iteration #6: *Charm!*

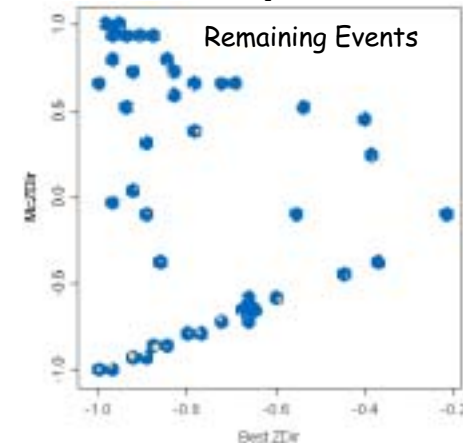1. Switch over to Onboard Flight Software Filter for "pruning"

   *Look Ahead*:

   Refiltered Events using FSW Filter *MINUS* bit #17 ("No Tracks")

   Kills - 3% of AG sample

   (Leaves $A_{eff}$ ~ 8000 cm$^2$ (E > 3 GeV)
   and $A_{eff}$ x $\Delta\Omega$ = 1.9 m$^2$-str)

   Kills - 60% of BGE sample (Rate: .03 Hz)



Remaining Events

2. Run at least 5X more events!   In fact we should consider simply starting
   a regular MC production regime rather then the current "one-off" approach

3. Explore alternative Variable Selection
   Mechanisms.

# Conclusions

- Not there yet....

- CT/RT Technology  Promising

- Need to condense various choices into data set(s) suitable for public consumption!